



Artificial Software Diversification for WebAssembly

JAVIER CABRERA-ARTEAGA

Doctoral Thesis
Supervised by
Benoit Baudry and Martin Monperrus
Stockholm, Sweden, 2023

KTH Royal Institute of Technology
School of Electrical Engineering and Computer Science
Division of Software and Computer Systems

TRITA-EECS-AVL-2020:4
ISBN 100-

SE-10044 Stockholm
Sweden

Akademisk avhandling som med tillstånd av Kungl Tekniska högskolan framlägges
till offentlig granskning för avläggande av Teknologie doktorexamen i elektroteknik
i .

© Javier Cabrera-Arteaga , date

Tryck: Universitetsservice US AB

Abstract

[1]

Keywords: Lorem, Ipsum, Dolor, Sit, Amet

Sammanfattning

[1]

LIST OF PAPERS

1. ***Superoptimization of WebAssembly Bytecode***
Javier Cabrera-Arteaga, Shrinish Donde, Jian Gu, Orestis Floros, Lucas Satabin, Benoit Baudry, Martin Monperrus
Conference Companion of the 4th International Conference on Art, Science, and Engineering of Programming (Programming 2021), MoreVMs
<https://doi.org/10.1145/3397537.3397567>
2. ***CROW: Code Diversification for WebAssembly***
Javier Cabrera-Arteaga, Orestis Floros, Oscar Vera-Pérez, Benoit Baudry, Martin Monperrus
Network and Distributed System Security Symposium (NDSS 2021), MADWeb
<https://doi.org/10.14722/madweb.2021.23004>
3. ***Multi-Variant Execution at the Edge***
Javier Cabrera-Arteaga, Pierre Laperdrix, Martin Monperrus, Benoit Baudry
Conference on Computer and Communications Security (CCS 2022), Moving Target Defense (MTD)
<https://dl.acm.org/doi/abs/10.1145/3560828.3564007>
4. ***WebAssembly Diversification for Malware Evasion***
Javier Cabrera-Arteaga, Tim Toady, Martin Monperrus, Benoit Baudry
Computers & Security, Volume 131, 2023
<https://www.sciencedirect.com/science/article/pii/S0167404823002067>
5. ***Wasm-mutate: Fast and Effective Binary Diversification for WebAssembly***
Javier Cabrera-Arteaga, Nick Fitzgerald, Martin Monperrus, Benoit Baudry
6. ***Scalable Comparison of JavaScript V8 Bytecode Traces***
Javier Cabrera-Arteaga, Martin Monperrus, Benoit Baudry
11th ACM SIGPLAN International Workshop on Virtual Machines and Intermediate Languages (SPLASH 2019)
<https://doi.org/10.1145/3358504.3361228>

ACKNOWLEDGEMENT

ACRONYMS

List of commonly used acronyms:

Wasm WebAssembly

Contents

List of Papers	iii
Acknowledgement	iv
Acronyms	v
Contents	1
I Thesis	2
1 Introduction	3
1.1 Background	3
1.2 Problem statement	3
1.3 Automatic Software diversification requirements	3
1.4 List of contributions	3
1.5 Summary of research papers	4
1.6 Thesis outline	4
2 Background and state of the art	5
2.1 WebAssembly	5
2.2 Software diversification	5
2.3 Generating Software Diversification	5
2.4 Exploiting Software Diversification	5
3 Automatic Software Diversification for WebAssembly	6
3.1 Compiler based approaches	6
3.2 MEWE: Multi-variant Execution for WebAssembly	12
3.3 Binary based approach	15

4	Evaluation	16
4.1	Use cases	16
4.2	Experimental protocols	16
4.3	Results	16
5	Results and discussion	17
5.1	Summary of technical contributions	17
5.2	Summary of empirical findings.	17
5.3	Summary of empirical findings.	17
5.4	Future Work	17
II	Included papers	18
	Superoptimization of WebAssembly Bytecode	20
	CROW: Code Diversification for WebAssembly	21
	Multi-Variant Execution at the Edge	22
	WebAssembly Diversification for Malware Evasion	23
	Wasm-mutate: Fast and Effective Binary Diversification for WebAssembly	24
	Scalable Comparison of JavaScript V8 Bytecode Traces	25

Part I

Thesis

TODO Recent papers first. Mention Workshops instead in conference. "Proceedings of XXXX". Add the pages in the papers list.

■ 1.1 Background

TODO Motivate with the open challenges.

■ 1.2 Problem statement

TODO Problem statement **TODO** Set the requirements as R1, R2, then map each contribution to them.

■ 1.3 Automatic Software diversification requirements

1. 1: **TODO** Requirement 1

■ 1.4 List of contributions

C1: Methodology contribution: We propose a methodology for generating software diversification for WebAssembly and the assessment of the generated diversity.

C2: Theoretical contribution: We propose theoretical foundation in order to improve Software Diversification for WebAssembly.

C3: Automatic diversity generation for WebAssembly: We generate WebAssembly program variants.

C4: Software Diversity for Defensive Purposes: We assess how generated WebAssembly program variants could be used for defensive purposes.

C5: Software Diversity for Offensives Purposes: We assess how generated WebAssembly program variants could be used for offensive purposes, yet improving security systems.

Contribution	Resarch papers				
	P1	P2	P3	P4	P5
C1	x	x		x	x
C2	x	x			
C3	x	x	x		
C4	x	x	x		
C5			x		
C6	x	x	x	x	x

Table 1.1: Mapping of the contributions to the research papers appended to this thesis.

C6: Software Artifacts: We provide software artifacts for the research community to reproduce our results.

TODO Make multi column table

■ 1.5 Summary of research papers

- P1:** Superoptimization of WebAssembly Bytecode.
- P2:** CROW: Code randomization for WebAssembly bytecode.
- P3:** Multivariant execution at the Edge.
- P4:** Wasm-mutate: Fast and efficient software diversification for WebAssembly.
- P5:** WebAssembly Diversification for Malware evasion.

■ 1.6 Thesis outline

- 2.1 WebAssembly

- 2.1.1 WebAssembly toolchains

- TODO** Mention, stress the landscape of tools that involve Wasm. Include analysis tools, fuzzers, optimizers and malware detectors.

- TODO** End up motivating the need of Software Diversification for: testing and reliability.

- 2.2 Software diversification

- 2.3 Generating Software Diversification

- 2.3.1 Variants generation

- 2.3.2 Variants equivalence

- 2.4 Exploiting Software Diversification

- 2.4.1 Defensive Diversification

- 2.4.2 Offensive Diversification

TODO Start here. 4 pages each and 2 pages discussion. Target 20 pages.

The work of Hilbig et al. [?] in 2021 influences our design decisions. According to their work, 70% of the Wasm binaries in the wild are created with LLVM-based compilers. Therefore, we provide artificial software diversity for Wasm through LLVM. Other solutions would have been to diversify at the source-code level or the Wasm binary level. However, these facts would limit the applicability of our work. Our approach is more general as diversification also will work for other LLVM backends.

LLVM is a compound of three main components [?]. First, the frontend (compilers such as clang and rustc) converts the program source code to LLVM intermediate representation (LLVM IR). Second, optimization and transformation processes improve the LLVM IR. Third and final, the backend component is in charge of generating the target machine code. In Figure 3.1 we show how we use the LLVM pipeline in our contributions, which are highlighted as dashed squares.

The global workflow in Figure 3.1 starts by receiving the source code. Then the LLVM frontend transforms it into LLVM IR representation ①. We alter the LLVM pipeline that compiles source code to Wasm by introducing a diversifier component.

The diversifier generates LLVM IR variants from the output of the frontend ②. The LLVM IR variants are inputs for our customized Wasm backend. The diversifier and the custom Wasm LLVM backend compose CROW, which creates Wasm program variants out of a source code program ③. In addition, an orthogonal tool comes from the generation of LLVM IR variants at Step ②. MEWE [?], merges and creates multivariant binaries to provide MVE for Wasm ④.

■ 3.1 Compiler based approaches

■ 3.1.3 CROW: Code Randomization of WebAssembly

This section describes the red squared tooling in Figure 3.1 named CROW [?]. CROW is a tool tailored to create semantically equivalent Wasm variants from

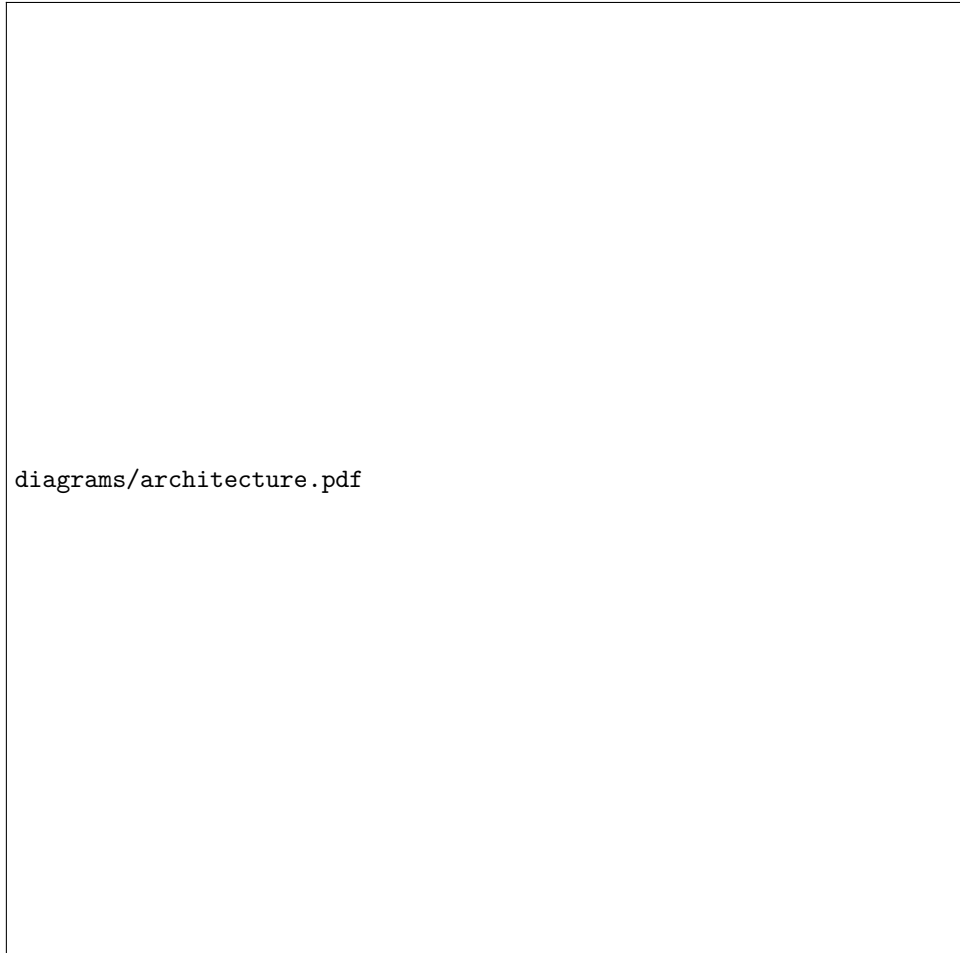


Figure 3.1: Generic workflow to create Wasm program variants.

an LLVM front-end output. Using a custom Wasm LLVM backend, it generates the Wasm binary variants.

In Figure 3.3, we describe the workflow of CROW to create program variants. The Diversifier in CROW is composed by two main processes, *exploration* and *combining*. The *exploration* process operates at the instruction level for each function in its input LLVM. For all LLVM instructions, CROW produces a collection of functionally equivalent code replacements. In the *combining* stage, CROW assembles the code replacements to generate different LLVM IR variants. CROW generates the LLVM IR variants by traversing the power set of all possible combinations of code replacements. Finally, the custom Wasm LLVM backend

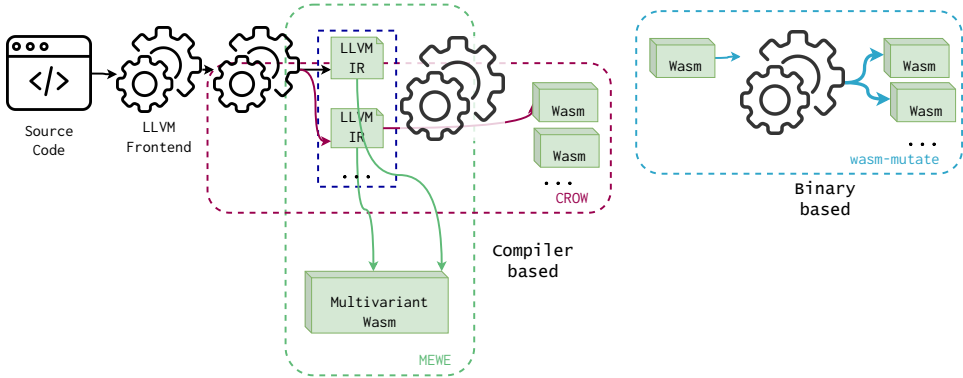


Figure 3.2: Approach landscape.

compiles the assembled LLVM IR variants into Wasm binaries. In the following text, we describe our design decisions. All our implementation choices are based on one premise: *each design decision should increase the number of Wasm variants that CROW creates.*

■ 3.1.4 Exploration

The primary component of CROW’s exploration process is its code replacements generation strategy. The diversifier implemented in CROW is based on the proposed superdiversifier of Jacob et al. [?]. A superoptimizer focuses on *searching* for a new program that is faster or smaller than the original code while preserving its functionality. The concept of superoptimizing a program dates back to 1987, with the seminal work of Massalin [?] which proposes an exhaustive exploration of the solution space. The search space is defined by choosing a subset of the machine’s instruction set and generating combinations of optimized programs, sorted by code size in ascending order. If any of these programs is found to perform the same function as the source program, the search halts. On the contrary, a superdiversifier keeps all intermediate search results despite their performance.

We use the superdiversifier idea of Jacob and colleagues to implement CROW because of two main reasons. First, the code replacements generated by this technique outperform diversification strategies based on handwritten rules. Concretely, we can control the quality of the generated codes. Besides, CROW always generates equivalent programs because it is based on a solver to check for equivalence. Second, there is a battle-tested superoptimizer for LLVM, Souper [?]. This latter makes it feasible the construction of a generic LLVM superdiversifier.



Figure 3.3: CROW components following the diagram in Figure 3.1. CROW takes LLVM IR to generate functionally equivalent code replacements. Then, CROW assembles program variants by combining them.

We modify Souper to keep all possible solutions in their searching algorithm. Souper builds a Data Flow Graph for each LLVM integer-returning instruction. Then, for each Data Flow Graph, Souper exhaustively builds all possible expressions from a subset of the LLVM IR language. Each syntactically correct expression in the search space is semantically checked versus the original with a theorem solver. Souper synthesizes the replacements in increasing size. Thus, the first found equivalent transformation is the optimal replacement result of the searching. CROW keeps more equivalent replacements during the searching by

removing the halting criteria. Instead the original halting conditions, CROW does not halt when it finds the first replacement. CROW continues the search until a timeout is reached or the replacements grow to a size larger than a predefined threshold.

Notice that the searching space increases exponentially with the size of the LLVM IR language subset. Thus, we prevent Souper from synthesizing instructions with no correspondence in the Wasm backend. This decision reduces the searching space. For example, creating an expression having the **freeze** LLVM instructions will increase the searching space for instruction without a Wasm's opcode in the end. Moreover, we disable the majority of the pruning strategies of Souper for the sake of more program variants. For example, Souper prevents the generation of the commutative operations during the searching. On the contrary, CROW still uses such transformation as a strategy to generate program variants.

■ 3.1.4 Constant inferring

One of the code transformation strategies of Souper does *constant inferring*. This means that Souper infers pieces of code as a single constant assignment. In particular, Souper focuses on variables that are used to control branches. By extending Souper as a superdiversifier, we add this transformation strategy as a new mutation strategy to the ones defined in ??.

After a *constant inferring*, the generated program is considerably different from the original program, being suitable for diversification. Let us illustrate the case with an example. The Babbage problem code in Listing 3.1 is composed of a loop that stops when it discovers the smaller number that fits with the Babbage condition in Line 4.

Listing 3.1: Babbage problem.

```

1      int babbage() {
2          int current = 0,
3              square;
4          while ((square=current*current) %
               ↪ 1000000 != 269696) {
5              current++;
6          }
7          printf ("The number is %d\n", current)
               ↪ ;
8          return 0 ;
9      }
```

In theory, this value can also be inferred

Listing 3.2: Constant inferring transformation over the original Babbage problem in Listing 3.1.

```
int babbage() {
    int current = 25264;

    printf ("The number is %d\n", current);
    return 0 ;
}
```

by unrolling the loop the correct number of times with the LLVM toolchain. However, standard LLVM tools cannot unroll the **while**-loop because the loop count is too large. The original Souper deals with this case, generating the program in Listing 3.2. It infers the value of **current** in Line 2 such that the Babbage condition is reached. Therefore, the condition in the loop will always be false. Then, the loop is dead code and is removed in the final compilation. The new program in Listing 3.2 is remarkably smaller and faster than the original code. Therefore, it offers differences both statically and at runtime¹.

■ 3.1.4 Removing subsequent optimizations for LLVM

During the implementation of CROW, we have the premise of removing all built-in optimizations in the LLVM backend that could reverse Wasm variants. Therefore, we modify the Wasm backend. We disable all optimizations in the Wasm backend that could reverse the CROW transformations. In the following enumeration, we list three concrete optimizations that we remove from the Wasm backend.²

- Constant folding: this optimization calculates the operation over two (or more) constants in compiling time, and replaces the original expression by its constant result. For example, let us suppose $a = 10 + 12$ a subexpression to be compiled, with the original optimization, the Wasm backend replaces it by $a = 22$.
- Expressions normalization: in this case, the comparison operations are normalized to its complementary operation, e.g. $a > b$ is always replaced by $b \leq a$.
- Redundant operation removal: expressions such as the multiplication of variables by $a = b2^n$ are replaced by shift left operations $a = b \ll n$.

¹Notice that for the sake of illustration, we show both codes in C language, this process inside CROW is performed directly in LLVM IR. Also, notice that the two programs in the example follow the definition of *functional equivalence* discussed in ??.

²We only illustrate three of the removed optimization for the sake of simplicity.

■ 3.2 MEWE: Multi-variant Execution for WebAssembly

This section describes MEWE [?]. MEWE synthesizes diversified function variants by using CROW. It then provides execution-path randomization in a Multivariant Execution (MVE). The tool generates application-level multivariant binaries without changing the operating system or Wasm runtime. MEWE creates an MVE by intermixing functions for which CROW generates variants, as step ② in Figure 3.1 shows. CROW generates each one of these variants with fine-grained diversification at the instruction level, applying the majority of the strategies discussed in ?? and *constant inferring*. MEWE adds a new mutation strategy. It inlines function variants when appropriate, resulting in call stack diversification at runtime.

In Figure 3.4 we zoom MEWE from the blue highlighted square in Figure 3.1. MEWE takes the LLVM IR variants generated by CROW’s diversifier. It then merges LLVM IR variants into a Wasm multivariant. In the figure, we highlight the two components of MEWE, *Multivariant Generation* and the *Mixer*. In the *Multivariant Generation* process, MEWE merges the LLVM IR variants created by CROW and creates an LLVM multivariant binary. The merging of the variants intermixes the calling of function variants, allowing the execution path randomization.

The Mixer augments the LLVM multivariant binary with a random generator. The random generator is needed to perform the execution-path randomization. Also, *The Mixer* fixes the entrypoint in the multivariant binary. Finally, MEWE generates a standalone multivariant Wasm binary using the same custom Wasm LLVM backend from CROW. Once generated, the multivariant Wasm binary can be deployed to any Wasm engine.

■ 3.2.1 Multivariant generation

The key component of MEWE consists in combining the variants into a single binary. The goal is to support execution-path randomization at runtime. The core idea is to introduce one dispatcher function per original function with variants. A dispatcher function is a synthetic function in charge of choosing a variant at random when the original function is called. With the introduction of the dispatcher function, MEWE turns the original call graph into a multivariant call graph, defined as follows.

Definition 1. *Multivariant Call Graph (MCG): A multivariant call graph is a call graph $\langle N, E \rangle$ where the nodes in N represent all the functions in the binary and an edge $(f_1, f_2) \in E$ represents a possible invocation of f_2 by f_1 [?]. The nodes in N have three possible types: a function present in the original program, a generated function variant, or a dispatcher function.*

In Figure 3.5, we show the original static call graph for an original program (top of the figure), as well as the multivariant call graph generated with MEWE

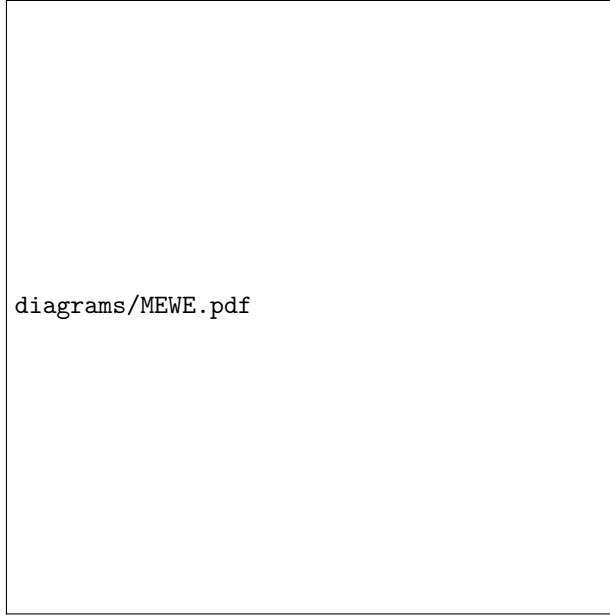


Figure 3.4: Overview of MEWE workflow. It takes as input an LLVM binary. It first generates a set of functionally equivalent variants for each function in the binary using CROW. Then, MEWE generates an LLVM multivariant binary composed of all the function variants. Finally, the Mixer includes the behavior in charge of selecting a variant when a function is invoked. Finally, the MEWE mixer composes the LLVM multivariant binary with a random number generation library and tampers the original application entrypoint. The final process produces a Wasm multivariant binary ready to be deployed.

(bottom of the figure). The gray nodes represent function variants, the green nodes function dispatchers, and the yellow nodes are the original functions. The directed edges represent the possible calls. The original program includes three functions. MEWE generates 43 variants for the first function, none for the second, and three for the third. MEWE introduces two dispatcher nodes for the first and third functions. Each dispatcher is connected to the corresponding function variants to invoke one variant randomly at runtime.

In Listing 3.3, we illustrate the LLVM construction for the function dispatcher corresponding to the right most green node of Figure 3.5. It first calls the random generator, which returns a value used to invoke a specific function variant. We implement the dispatchers with a switch-case structure to avoid indirect calls that can be susceptible to speculative execution-based attacks [?]. The choice of a switch-case also avoids having multiple function definitions with the same signature, which could increase the attack surface in case the function signature

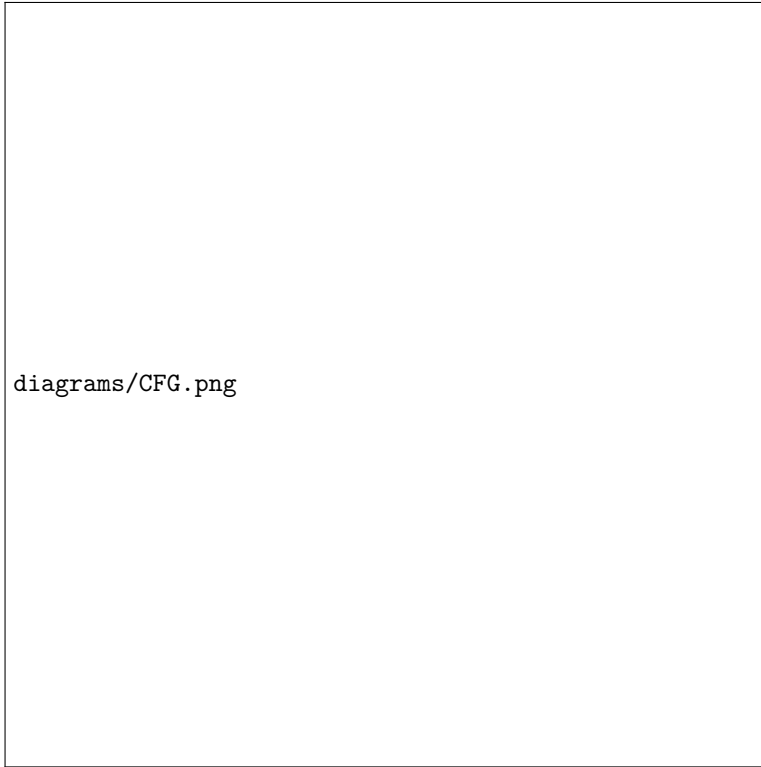


Figure 3.5: Example of two static call graphs. At the top, the original call graph, at the bottom, the multivariant call graph, which includes nodes that represent function variants (in gray), dispatchers (in green), and original functions (in yellow).

is vulnerable [?]. This also allows MEWE to inline function variants inside the dispatcher instead of defining them again. Here we trade security over performance since dispatcher functions that perform indirect calls, instead of a switch-case, could improve the performance of the dispatchers as indirect calls have constant time.

■ 3.2.2 The Mixer

MEWE has four specific objectives: link the LLVM multivariant binary, inject a random generator, tamper the application’s entrypoint, and merge all these components into a multivariant Wasm binary. We use the Rustc compiler³ to orchestrate the mixing. For the random generator, we rely on WASI’s specification [?] for the random behavior of the dispatchers. However, its exact

³<https://doc.rust-lang.org/rustc/what-is-rustc.html>

```
define internal i32 @foo(i32 %0) {
  entry:
    %1 = call i32 @discriminate(i32 3)
    switch i32 %1, label %end [
      i32 0, label %case_43_
      i32 1, label %case_44_
    ]
  case_43_:
    %2 = call i32 @foo_43_(%0)
    ret i32 %2
  case_44_:
    %3 = <body of foo_44_ inlined>
    ret i32 %3
end:
  %4 = call i32 @foo_original(%0)
  ret i32 %4
}
```

Listing 3.3: Dispatcher function embedded in the multivariant binary of the original function in the rightmost green node in Figure 3.5.

implementation is dependent on the platform on which the binary is deployed. The Mixer creates a new entrypoint for the binary called *entrypoint tampering*. It wraps the dispatcher for the entrypoint variants as a new function for the final Wasm binary and is declared as the application entrypoint.

■ 3.3 Binary based approach

■ 3.3.1 wasm-mutate

■ 4.1 Use cases

RQ1: Defensive Diversification: ?

RQ2: Offensive Diversification: ?

■ 4.2 Experimental protocols

■ 4.2.0 Metrics

New static metric. Diversification preservation.

■ 4.3 Results

- 5.1 Summary of technical contributions
- 5.2 Summary of empirical findings
- 5.3 Summary of empirical findings
- 5.4 Future Work

REFERENCES

Part II

Included papers

SUPEROPTIMIZATION OF WEBASSEMBLY BYTECODE

Javier Cabrera-Arteaga, Shrinish Donde, Jian Gu, Orestis Floros, Lucas Satabin, Benoit Baudry, Martin Monperrus

Conference Companion of the 4th International Conference on Art, Science, and Engineering of Programming (Programming 2021), MoreVMs

<https://doi.org/10.1145/3397537.3397567>

CROW: CODE DIVERSIFICATION FOR WEBASSEMBLY

Javier Cabrera-Arteaga, Orestis Floros, Oscar Vera-Pérez, Benoit Baudry,
Martin Monperrus

Network and Distributed System Security Symposium (NDSS 2021), MADWeb

<https://doi.org/10.14722/madweb.2021.23004>

MULTI-VARIANT EXECUTION AT THE EDGE

Javier Cabrera-Arteaga, Pierre Laperdrix, Martin Monperrus, Benoit Baudry
*Conference on Computer and Communications Security (CCS 2022), Moving
Target Defense (MTD)*

<https://dl.acm.org/doi/abs/10.1145/3560828.3564007>

WEBASSEMBLY DIVERSIFICATION FOR MALWARE EVASION

Javier Cabrera-Arteaga, Tim Toady, Martin Monperrus, Benoit Baudry
Computers & Security, Volume 131, 2023

<https://www.sciencedirect.com/science/article/pii/S0167404823002067>

WASM-MUTATE: FAST AND EFFECTIVE BINARY DIVERSIFICATION FOR WEBASSEMBLY

Javier Cabrera-Arteaga, Nick Fitzgerald, Martin Monperrus, Benoit Baudry
Under revision

SCALABLE COMPARISON OF JAVASCRIPT V8 BYTECODE TRACES

Javier Cabrera-Arteaga, Martin Monperrus, Benoit Baudry
*11th ACM SIGPLAN International Workshop on Virtual Machines and
Intermediate Languages (SPLASH 2019)*

<https://doi.org/10.1145/3358504.3361228>